

Nonparametric Reduced Rank Regression

Rina Foygel^{†,*}, Michael Horrell[†], Mathias Drton^{†,‡} and John Lafferty^{†^B}

^{*}*Department of Statistics
Stanford University*

[†]*Department of Statistics
^BDepartment of Computer Science
The University of Chicago*

[‡]*Department of Statistics
University of Washington*

January 10, 2013

Abstract: We propose an approach to multivariate nonparametric regression that generalizes reduced rank regression for linear models. An additive model is estimated for each dimension of a q -dimensional response, with a shared p -dimensional predictor variable. To control the complexity of the model, we employ a functional form of the Ky-Fan or nuclear norm, resulting in a set of function estimates that have low rank. Backfitting algorithms are derived and justified using a nonparametric form of the nuclear norm subdifferential. Oracle inequalities on excess risk are derived that exhibit the scaling behavior of the procedure in the high dimensional setting. The methods are illustrated on gene expression data.

Keywords and phrases: multivariate regression, nonparametric regression, nuclear norm regularization, high dimensional inference, oracle risk bounds.

1. Introduction

In the multivariate regression problem the objective is to estimate the conditional mean

$$\mathbb{E}(Y | X) = m(X) = (m^1(X), \dots, m^q(X))^\top,$$

where Y is a q -dimensional response vector, X is a p -dimensional covariate vector, and we are given a sample of n i.i.d. pairs $\{(X^{(i)}, Y^{(i)})\}$ from the joint distribution of X and Y . This is also referred to as multi-task learning in the machine learning literature. Under a linear model, the mean is estimated as $m(X) = BX$ where $B \in \mathbb{R}^{q \times p}$ is a $q \times p$ matrix of regression coefficients. When the dimensions p and q are large relative to the sample size n , the coefficients of B cannot be reliably estimated, without further assumptions.

In reduced rank regression the matrix B is estimated under a rank constraint $r = \text{rank}(B) \leq C$, so that the rows or columns of B lie in an r -dimensional subspace of \mathbb{R}^q or \mathbb{R}^p . Intuitively, this implies that the model is based on a smaller number of features than

This is an extended version of a paper presented at NIPS (Foygel *et al.*, 2012).

the ambient dimensionality p would suggest, or that the tasks representing the components Y^k of the response are closely related. In low dimensions, the constrained rank model can be computed as an orthogonal projection of the least squares solution; but in high dimensions this is not well defined.

Recent research has studied the use of the nuclear norm as a convex surrogate for the rank constraint. The nuclear norm $\|B\|_*$, also known as the trace or Ky-Fan norm, is the sum of the singular vectors of B . A rank constraint can be thought of as imposing sparsity, but in an unknown basis; the nuclear norm plays the role of the ℓ_1 norm in sparse estimation. Its use for low rank estimation problems was proposed by [Fazel \(2002\)](#). More recently, nuclear norm regularization in multivariate linear regression has been studied by [Yuan *et al.* \(2007\)](#), and by [Negahban and Wainwright \(2011\)](#), who analyzed the scaling properties of the procedure in high dimensions.

In this paper we study nonparametric parallels of reduced rank linear models. We focus our attention on additive models, so that the regression function $m(X) = (m^1(X), \dots, m^q(X))^T$ has each component $m^k(X) = \sum_{j=1}^p m_j^k(X_j)$ equal to a sum of p functions, one for each covariate. The objective is then to estimate the $q \times p$ matrix of functions $M(X) = [m_j^k(X_j)]$.

The first problem we address, in [Section 2](#), is to determine a replacement for the regularization penalty $\|B\|_*$ in the linear model. Because we must estimate a matrix of functions, the analogue of the nuclear norm is not immediately apparent. We propose two related regularization penalties for nonparametric low rank regression, and show how they specialize to the linear case. We then study, in [Section 4](#), the (infinite dimensional) subdifferential of these penalties. In the population setting, this leads to stationary conditions for the minimizer of the regularized mean squared error. This subdifferential calculus then justifies penalized backfitting algorithms for carrying out the optimization for a finite sample. Constrained rank additive models (CRAM) for multivariate regression are analogous to sparse additive models (SPAM) for the case where the response is 1-dimensional ([Ravikumar *et al.*, 2009](#)) (studied also in the reproducing kernel Hilbert space setting by [Raskutti, Wainwright and Yu \(2012\)](#)), but with the goal of recovering a low-rank matrix rather than an entry-wise sparse vector. The backfitting algorithms we derive in [Section 5](#) are analogous to the iterative smoothing and soft thresholding backfitting algorithms for SPAM proposed by [Ravikumar *et al.* \(2009\)](#). A uniform bound on the excess risk of the estimator relative to an oracle is given [Section 7](#). This shows the statistical scaling behavior of the methods for prediction. The analysis requires a concentration result for nonparametric covariance matrices in the spectral norm. Experiments with synthetic, gene, and biochemistry data are given in [Section 8](#), which are used to illustrate different facets of the proposed nonparametric reduced rank regression techniques.

2. Nonparametric Nuclear Norm Penalization

We begin by presenting the penalty that we will use to induce nonparametric regression estimates to be low rank. To motivate our choice of penalty and provide some intuition, suppose that $f^1(x), \dots, f^q(x)$ are q smooth one-dimensional functions with a common domain. What does it mean for this collection of functions to be low rank? Let $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ be a collection of points in the common domain of the functions. We require that the $n \times q$ matrix of function values $\mathbb{F}(x^{1:n}) = [f^k(x^{(i)})]$ is low rank. This matrix is of rank at most $r < q$ for every set $\{x^{(i)}\}$ of arbitrary size n if and only if the functions $\{f^k\}$ are r -linearly independent—each function can be written as a linear combination of r of the other functions.

In the multivariate regression setting, but still assuming the domain is one-dimensional for simplicity ($q > 1$ and $p = 1$), we have a random sample $X^{(1)}, \dots, X^{(n)}$. Consider the $n \times q$ sample matrix $\mathbb{M} = [m^k(X^{(i)})]$ associated with a vector $M = (m^1, \dots, m^q)$ of q smooth (regression) functions, and suppose that $n > q$. We would like for this to be a low rank matrix. This suggests the penalty

$$\|\mathbb{M}\|_* = \sum_{s=1}^q \sigma_s(\mathbb{M}) = \sum_{s=1}^q \sqrt{\lambda_s(\mathbb{M}^\top \mathbb{M})},$$

where $\{\lambda_s(A)\}$ denotes the eigenvalues of a symmetric matrix A and $\{\sigma_s(B)\}$ denotes the singular values of a matrix B . Now, assuming the columns of \mathbb{M} are centered, and $\mathbb{E}[m^k(X)] = 0$ for each k , we recognize $\frac{1}{n}\mathbb{M}^\top \mathbb{M}$ as the sample covariance $\widehat{\Sigma}(M)$ of the population covariance

$$\Sigma(M) := \text{Cov}(M(X)) = [\mathbb{E}(m^k(X)m^l(X))].$$

This motivates the following sample and population penalties, where $A^{1/2}$ denotes the matrix square root:

$$\text{population penalty: } \|\Sigma(M)^{1/2}\|_* = \|\text{Cov}(M(X))^{1/2}\|_* \quad (2.1)$$

$$\text{sample penalty: } \|\widehat{\Sigma}(M)^{1/2}\|_* = \frac{1}{\sqrt{n}}\|\mathbb{M}\|_*. \quad (2.2)$$

We will also use the notation $\|M\|_* = \|\text{Cov}(M(X))^{1/2}\|_*$.

This leads to the following population and empirical regularized risk functionals for low rank nonparametric regression:

$$\text{population penalized risk: } \frac{1}{2}\mathbb{E}\|Y - M(X)\|_2^2 + \lambda\|\Sigma(M)^{1/2}\|_* \quad (2.3)$$

$$\text{empirical penalized risk: } \frac{1}{2n}\|Y - \mathbb{M}\|_F^2 + \frac{\lambda}{\sqrt{n}}\|\mathbb{M}\|_*, \quad (2.4)$$

where, in the empirical case, Y denotes the $n \times q$ matrix of response values for the sample $\{(X^{(i)}, Y^{(i)})\}$. We recall that if $A \succeq 0$ has spectral decomposition $A = UDU^\top$ then $A^{1/2} = UD^{1/2}U^\top$.

3. Constrained Rank Additive Models (CRAM)

We now consider the case where X is p -dimensional. Throughout the paper we use superscripts to denote indices of the q -dimensional response, and subscripts to denote indices of the p -dimensional covariate. We consider the family of additive models, with regression functions of the form

$$m(X) = (m^1(X), \dots, m^q(X))^\top = \sum_{j=1}^p M_j(X_j),$$

each term $M_j(X_j) = (m_j^1(X_j), \dots, m_j^q(X_j))^\top$ being a q -vector of functions evaluated at X_j .

In this setting we propose two different penalties. The first penalty, intuitively, encourages the vector $(m_j^1(X_j), \dots, m_j^q(X_j))$ to be low rank, for each j . Assume that the functions $m_j^k(X_j)$ all have mean zero; this is required for identifiability in the additive model. As a shorthand, let $\Sigma_j = \Sigma(M_j) = \text{Cov}(M_j(X_j))$ denote the covariance matrix of the j -th component functions, with sample version $\hat{\Sigma}_j$. The population and sample versions of the first penalty are then given by

$$\|\Sigma_1^{1/2}\|_* + \|\Sigma_2^{1/2}\|_* + \dots + \|\Sigma_p^{1/2}\|_* \quad (3.1)$$

$$\|\hat{\Sigma}_1^{1/2}\|_* + \|\hat{\Sigma}_2^{1/2}\|_* + \dots + \|\hat{\Sigma}_p^{1/2}\|_* = \frac{1}{\sqrt{n}} \sum_{j=1}^p \|\mathbb{M}_j\|_*. \quad (3.2)$$

The second penalty encourages the set of q vector-valued functions $(m_1^k, m_2^k, \dots, m_p^k)^\top$ to be low rank. This penalty is given by

$$\left\| \left(\Sigma_1^{1/2} \dots \Sigma_p^{1/2} \right) \right\|_* \quad (3.3)$$

$$\left\| \left(\hat{\Sigma}_1^{1/2} \dots \hat{\Sigma}_p^{1/2} \right) \right\|_* = \frac{1}{\sqrt{n}} \|\mathbb{M}_{1:p}\|_* \quad (3.4)$$

where, for convenience of notation, $\mathbb{M}_{1:p} = (\mathbb{M}_1^\top \dots \mathbb{M}_p^\top)^\top$ is an $np \times q$ matrix. The corresponding population and empirical risk functionals, for the first penalty, are then

$$\frac{1}{2} \mathbb{E} \left\| Y - \sum_{j=1}^p M_j(X) \right\|_2^2 + \lambda \sum_{j=1}^p \|\Sigma_j^{1/2}\|_* \quad (3.5)$$

$$\frac{1}{2n} \left\| Y - \sum_{j=1}^p \mathbb{M}_j \right\|_F^2 + \frac{\lambda}{\sqrt{n}} \sum_{j=1}^p \|\mathbb{M}_j\|_* \quad (3.6)$$

and similarly for the second penalty.

Now suppose that each X_j is normalized so that $\mathbb{E}(X_j^2) = 1$. In the linear case we have $M_j(X_j) = X_j B_j$ where $B_j \in \mathbb{R}^q$. Let $B = (B_1 \dots B_p) \in \mathbb{R}^{q \times p}$. Some straightforward calculation shows that the penalties reduce to

$$\|\Sigma_j^{1/2}\|_* = \|B_j\|_2 \quad (3.7)$$

$$\|\Sigma_1^{1/2} \dots \Sigma_p^{1/2}\|_* = \|B\|_*. \quad (3.8)$$

Thus, in the linear case the first penalty is encouraging B to be column-wise sparse, so that many of the B_j s are zero, meaning that X_j doesn't appear in the fit. This is a version of the group lasso (Yuan and Lin, 2006). The second penalty reduces to the nuclear norm regularization $\|B\|_*$ used for high-dimensional reduced-rank regression.

4. Subdifferentials for Functional Matrix Norms

A key to deriving algorithms for functional low-rank regression is computation of the subdifferentials of the penalties. We are interested in $(q \times p)$ -dimensional matrices of functions $F = [f_j^k]$. For each column index j and row index k , f_j^k is a function of a random variable X_j , and we will take expectations with respect to X_j implicitly. We write F_j to mean the j th column of F , which is a q -vector of functions of X_j . We define the inner product between two matrices of functions as

$$\langle\langle F, G \rangle\rangle := \sum_{j=1}^p \sum_{k=1}^q \mathbb{E}(f_j^k g_j^k) = \sum_{j=1}^p \mathbb{E}(F_j^\top G_j) = \text{tr}(\mathbb{E}(F G^\top)) , \quad (4.1)$$

and write $\|F\|_2 = \sqrt{\langle\langle F, F \rangle\rangle}$. Note that $\|F\|_2$ equals the Frobenius norm of $\sqrt{\mathbb{E}(F F^\top)}$ where $\mathbb{E}(F F^\top) = \sum_j \mathbb{E}(F_j F_j^\top) \succeq 0$ is a positive semidefinite $q \times q$ matrix.

We define two further norms on a matrix of functions F , namely,

$$\|F\|_{\text{sp}} := \sqrt{\|\mathbb{E}(F F^\top)\|_{\text{sp}}} = \left\| \sqrt{\mathbb{E}(F F^\top)} \right\|_{\text{sp}} \quad \text{and} \quad \|F\|_* := \left\| \sqrt{\mathbb{E}(F F^\top)} \right\|_*,$$

where $\|A\|_{\text{sp}}$ is the spectral norm (operator norm), the largest singular value of A , and it is convenient to write the matrix square root as $\sqrt{A} = A^{1/2}$. Each of the norms depends on F only through $\mathbb{E}(F F^\top)$. In fact, these two norms are dual—for any F ,

$$\|F\|_* = \sup_{\|G\|_{\text{sp}} \leq 1} \langle\langle G, F \rangle\rangle , \quad (4.2)$$

where the supremum is attained by setting $G = \left(\sqrt{\mathbb{E}(F F^\top)} \right)^{-1} F$, with A^{-1} denoting the matrix pseudo-inverse.

Proposition 4.1. *The subdifferential of $\|F\|_*$ is the set*

$$\mathcal{S}(F) := \left\{ \left(\sqrt{\mathbb{E}(F F^\top)} \right)^{-1} F + H : \|H\|_{\text{sp}} \leq 1, \mathbb{E}(F H^\top) = \mathbf{0}_{q \times q}, \mathbb{E}(F F^\top) H = \mathbf{0}_{q \times p} \text{ a.e.} \right\} . \quad (4.3)$$

Proof. The fact that $\mathcal{S}(F)$ contains the subdifferential $\partial\|F\|_*$ can be proved by comparing our setting (matrices of functions) to the ordinary matrix case; see Recht, Fazel and Parrilo (2010); Watson (1992). Here, we show the reverse inclusion, $\mathcal{S}(F) \subseteq \partial\|F\|_*$. Let $D \in \mathcal{S}(F)$ and let G be any element of the function space. We need to show

$$\|F + G\|_* \geq \|F\|_* + \langle\langle G, D \rangle\rangle , \quad (4.4)$$

where $D = \left(\sqrt{\mathbb{E}(FF^\top)}\right)^{-1} F + H =: \tilde{F} + H$ for some H satisfying the conditions in (4.3) above. Expanding the right-hand side of (4.4), we have

$$\|F\|_* + \langle\langle G, D \rangle\rangle = \|F\|_* + \langle\langle G, \tilde{F} + H \rangle\rangle \quad (4.5)$$

$$= \langle\langle F + G, \tilde{F} + H \rangle\rangle \quad (4.6)$$

$$\leq \|F + G\|_* \|D\|_{\text{sp}}, \quad (4.7)$$

where the second equality follows from $\|F\|_* = \langle\langle F, \tilde{F} \rangle\rangle$, and the fact that $\langle\langle F, H \rangle\rangle = \text{tr}(\mathbb{E}(FH^\top)) = 0$. The inequality follows from the duality of the norms.

Finally, we show that $\|D\|_{\text{sp}} \leq 1$. We have

$$\mathbb{E}(DD^\top) = \mathbb{E}(\tilde{F}\tilde{F}^\top) + \mathbb{E}(\tilde{F}H^\top) + \mathbb{E}(H\tilde{F}^\top) + \mathbb{E}(HH^\top) \quad (4.8)$$

$$= \mathbb{E}(\tilde{F}\tilde{F}^\top) + \mathbb{E}(HH^\top), \quad (4.9)$$

where we use the fact that $\mathbb{E}(FH^\top) = \mathbf{0}_{q \times q}$, implying $\mathbb{E}(\tilde{F}H^\top) = \mathbf{0}_{q \times q}$. Next, let $\mathbb{E}(FF^\top) = VDV^\top$ be a reduced singular value decomposition, where D is a positive diagonal matrix of size $q' \leq q$. Then $\mathbb{E}(\tilde{F}\tilde{F}^\top) = VV^\top$, and we have

$$\mathbb{E}(FF^\top) \cdot H = \mathbf{0}_{q \times p} \text{ a.e.} \Leftrightarrow V^\top H = \mathbf{0}_{q' \times p} \text{ a.e.} \Leftrightarrow \mathbb{E}(\tilde{F}\tilde{F}^\top)H = \mathbf{0}_{q \times p} \text{ a.e.}.$$

This implies that $\mathbb{E}(\tilde{F}\tilde{F}^\top) \cdot \mathbb{E}(HH^\top) = \mathbf{0}_{q \times q}$ and so these two symmetric matrices have orthogonal row spans and orthogonal column spans. Therefore,

$$\|\mathbb{E}(DD^\top)\|_{\text{sp}} = \|\mathbb{E}(\tilde{F}\tilde{F}^\top) + \mathbb{E}(HH^\top)\|_{\text{sp}} \quad (4.10)$$

$$= \max \left\{ \|\mathbb{E}(\tilde{F}\tilde{F}^\top)\|_{\text{sp}}, \|\mathbb{E}(HH^\top)\|_{\text{sp}} \right\} \quad (4.11)$$

$$\leq 1, \quad (4.12)$$

where the last bound comes from the fact that $\|\tilde{F}\|_{\text{sp}}, \|H\|_{\text{sp}} \leq 1$. Therefore $\|D\|_{\text{sp}} \leq 1$. \square

This gives the subdifferential of penalty 2, defined in (3.3). We can view the first penalty update as just a special case of the second penalty update. For penalty 1 in (3.1), if we are updating F_j and fix all the other functions, we are now penalizing the norm

$$\|F_j\|_* = \left\| \sqrt{\mathbb{E}(F_j F_j^\top)} \right\|_* , \quad (4.13)$$

which is clearly just a special case of penalty 2 with a single q -vector of functions instead of p different q -vectors of functions. So, we have

$$\partial \|F_j\|_* = \left\{ \left(\sqrt{\mathbb{E}(F_j F_j^\top)} \right)^{-1} F_j + H_j : \|H_j\|_{\text{sp}} \leq 1, \mathbb{E}(F_j H_j^\top) = \mathbf{0}, \mathbb{E}(F_j F_j^\top) H_j = \mathbf{0} \text{ a.e.} \right\}.$$

5. Stationary Conditions and Backfitting Algorithms

Returning to the base case of $p = 1$ covariate, consider the population regularized risk optimization

$$\min_M \left\{ \frac{1}{2} \mathbb{E} \|Y - M(X)\|_2^2 + \lambda \|M\|_* \right\}, \quad (5.1)$$

where M is a vector of q univariate functions. The stationary condition for this optimization is

$$\mathbb{E}(Y | X) = M(X) + \lambda V(X) \quad \text{a.e. for some } V \in \partial \|M\|_*. \quad (5.2)$$

Define $P(X) := \mathbb{E}(Y | X)$.

Proposition 5.1. *Let $\mathbb{E}(PP^\top) = U \text{diag}(\tau) U^\top$ be the singular value decomposition and define*

$$M = U \text{diag}([1 - \lambda/\sqrt{\tau}]_+) U^\top P \quad (5.3)$$

where $[x]_+ = \max(x, 0)$. Then M satisfies stationary condition (5.2), and is a minimizer of (5.1).

Proof. Assume the singular values are sorted as $\tau_1 \geq \tau_2 \geq \dots \geq \tau_q$, and let r be the largest index such that $\sqrt{\tau_r} > \lambda$. Thus, M has rank r . Note that $\sqrt{\mathbb{E}(MM^\top)} = U \text{diag}([\sqrt{\tau} - \lambda]_+) U^\top$, and therefore

$$\lambda (\sqrt{\mathbb{E}(MM^\top)})^{-1} M = U \text{diag}(\lambda/\sqrt{\tau_{1:r}}, \mathbf{0}_{q-r}) U^\top P \quad (5.4)$$

where $x_{1:k} = (x_1, \dots, x_k)$ and $c_k = (c, \dots, c)$. It follows that

$$M + \lambda (\sqrt{\mathbb{E}(MM^\top)})^{-1} M = U \text{diag}(\mathbf{1}_r, \mathbf{0}_{q-r}) U^\top P. \quad (5.5)$$

Now define

$$H = \frac{1}{\lambda} U \text{diag}(\mathbf{0}_r, \mathbf{1}_{q-r}) U^\top P \quad (5.6)$$

and take $V = (\sqrt{\mathbb{E}(MM^\top)})^{-1} M + H$. Then we have $M + \lambda V = P$.

It remains to show that H satisfies the conditions of the subdifferential in (4.3). Since

$$\sqrt{\mathbb{E}(HH^\top)} = U \text{diag}(\mathbf{0}_r, \sqrt{\tau_{r+1}}/\lambda, \dots, \sqrt{\tau_q}/\lambda) U^\top \quad (5.7)$$

we have $\|H\|_{\text{sp}} \leq 1$. Also, $\mathbb{E}(MH^\top) = \mathbf{0}_{q \times q}$ since

$$\text{diag}(1 - \lambda/\sqrt{\tau_{1:r}}, \mathbf{0}_{q-r}) \text{diag}(\mathbf{0}_r, \mathbf{1}_{q-r}/\lambda) = \mathbf{0}_{q \times q}. \quad (5.8)$$

Similarly, $\mathbb{E}(MM^\top)H = \mathbf{0}_{q \times q}$ since

$$\text{diag}((\sqrt{\tau_{1:r}} - \lambda)^2, \mathbf{0}_{q-r}) \text{diag}(\mathbf{0}_r, \mathbf{1}_{q-r}/\lambda) = \mathbf{0}_{q \times q}. \quad (5.9)$$

It follows that $V \in \partial \|M\|_{\text{sp}}$. □

Input: Data matrices $X \in \mathbb{R}^{n \times p}$ and $Y \in \mathbb{R}^{n \times q}$, regularization parameter $\lambda \geq 0$.

Initialize $\widehat{\mathbb{M}}_j = (\widehat{m}_j^k(X_{ij})) = 0 \in \mathbb{R}^{n \times q}$, for $j = 1, \dots, p$.

Iterate until convergence:

For each $j = 1, \dots, p$:

- (1) Compute the residual $Z_j \leftarrow Y - \sum_{j' \neq j} \widehat{\mathbb{M}}_{j'}$.
- (2) Estimate $P_j = \mathbb{E}[Z_j | X_j]$ by smoothing: $\widehat{P}_j = \mathcal{S}_j Z_j$.
- (3) Compute SVD: $\frac{1}{n} \widehat{P}_j^\top \widehat{P}_j = U \text{diag}(\tau) U^\top$.
- (4) Soft threshold: $\widehat{\mathbb{M}}_j \leftarrow \widehat{P}_j U \text{diag}([1 - \lambda/\sqrt{\tau}]_+) U^\top$.
- (5) Center: $\widehat{\mathbb{M}}_j \leftarrow \widehat{\mathbb{M}}_j - \text{mean}(\widehat{\mathbb{M}}_j)$.

Output: Component functions $\widehat{\mathbb{M}}_j$ and estimates of the conditional mean vector $\sum_j \widehat{\mathbb{M}}_{ij}$.

FIG 1. The CRAM backfitting algorithm, using the first penalty, which penalizes each component.

The analysis above justifies a backfitting algorithm for estimating a constrained rank additive model with the first penalty, where the objective is

$$\min_{M_j} \left\{ \frac{1}{2} \mathbb{E} \left\| Y - \sum_{j=1}^p M_j(X_j) \right\|_2^2 + \lambda \sum_{j=1}^p \|M_j\|_* \right\}. \quad (5.10)$$

For a given coordinate j , we form the residual $Z_j = Y - \sum_{k \neq j} M_k$, and then compute the projection $P_j = \mathbb{E}(Z_j | X_j)$, with singular value decomposition $\mathbb{E}(P_j P_j^\top) = U \text{diag}(\tau) U^\top$. We then update

$$M_j = U \text{diag}([1 - \lambda/\sqrt{\tau}]_+) U^\top P_j \quad (5.11)$$

and proceed to the next variable. This is a Gauss-Seidel procedure that parallels the population backfitting algorithm for SPAM (Ravikumar *et al.*, 2009).

In the sample version we replace the conditional expectation $P_j = \mathbb{E}(Z_j | X_j)$ by a non-parametric linear smoother, $\widehat{P}_j = \mathcal{S}_j Z_j$. The algorithm is given in Figure 1. The algorithm for penalty 2 is similar and given in Figure 2. Both algorithms can be viewed as functional projected gradient descent procedures. Note that to predict at a point x not included in the training set, the smoother matrices are constructed using that point; that is, $\widehat{P}_j(x_j) = S_j(x_j)^\top Z_j$.

6. Working over an RKHS

Suppose that the functions m_j^k are required to lie in a Hilbert space \mathcal{H}_j . A modified empirical optimization is (for the first penalty)

$$\left\| Y - \sum_{j=1}^p \mathbb{M}_j \right\|_F^2 + \lambda_n \sum_{j=1}^p \|\mathbb{M}_j\|_* + \rho_n \sum_{j=1}^p \sum_{k=1}^q \|m_j^k\|_{\mathcal{H}_j} \quad (6.1)$$

Input: Data matrices $X \in \mathbb{R}^{n \times p}$ and $Y \in \mathbb{R}^{n \times q}$, regularization parameter $\lambda \geq 0$.

Initialize $\widehat{\mathbb{M}}_j = (\widehat{m}_j^k(X_{ij})) = 0 \in \mathbb{R}^{n \times q}$, for $j = 1, \dots, p$.

Iterate until convergence:

For each $j = 1, \dots, p$:

- (1) Compute the residual $Z_j \leftarrow Y - \sum_{j' \neq j} \widehat{\mathbb{M}}_{j'}$.
- (2) Estimate $P_j = \mathbb{E}[Z_j | X_j]$ by smoothing: $\widehat{P}_j = \mathcal{S}_j Z_j$.
- (3) Compute SVD: $\frac{1}{n} \widehat{P}_{1:p}^\top \widehat{P}_{1:p} = U \text{diag}(\tau) U^\top$.
- (4) Soft threshold: $\widehat{\mathbb{M}}_{1:p} \leftarrow \widehat{P}_{1:p} U \text{diag}([1 - \lambda/\sqrt{\tau}]_+) U^\top$.
- (5) Center: $\widehat{\mathbb{M}}_j \leftarrow \widehat{\mathbb{M}}_j - \text{mean}(\widehat{\mathbb{M}}_j)$.

Output: Component functions $\widehat{\mathbb{M}}_j$ and estimates of the conditional mean vector $\sum_j \widehat{\mathbb{M}}_{ij}$.

FIG 2. The CRAM backfitting algorithm, using the second penalty, which penalizes the components together.

where Y is an $n \times q$ data matrix and \mathbb{M}_j is an $n \times q$ matrix of function values associated with the j th columns of a $n \times p$ data matrix X . The first penalty is then a nuclear-norm constraint on these observed function values. The second penalty is a smoothness penalty on each of the coordinate functions in the appropriate Hilbert space, and is not empirical.

If \mathcal{H}_j is an RKHS with kernel K_j , then the representer theorem implies that we can restrict to functions m_j^k of the form

$$m_j^k(\cdot) = \sum_{i=1}^n \alpha_{ij}^k K_j(x_{ij}, \cdot), \quad (6.2)$$

where the α_{ij}^k are real weights. In this case the optimization becomes a finite dimensional semidefinite program over α . This parallels the approach of [Raskutti, Wainwright and Yu \(2012\)](#) for sparse additive models; see also [Dinuzzo and Fukumizu \(2012\)](#).

If $K_j = [K_{h_j}(x_{ij}, x_{i'j})] \in \mathbb{R}^{n \times n}$ denotes the Gram matrix for the j th variable, then $F_j = K_j \alpha_j$ where $\alpha_j = [\alpha_{ij}^k] \in \mathbb{R}^{n \times q}$. Using the first penalty, the convex optimization is then

$$\min_{\alpha} \frac{1}{2n} \left\| Y - \sum_j K_j \alpha_j \right\|_F^2 + \lambda_n \sum_{j=1}^p \|K_j \alpha_j\|_* + \rho_n \sum_{k=1}^q \sum_{j=1}^p \sqrt{\alpha_j^{kT} K_j \alpha_j^k} \quad (6.3)$$

where the third term is a smoothness penalty for the RKHS. This is a cone program with constraints involving both the second-order cone and the semidefinite cone.

7. Excess Risk Bounds

The population risk of a $q \times p$ regression matrix $M(X) = [M_1(X_1) \cdots M_p(X_p)]$ is

$$R(M) = \mathbb{E} \|Y - M(X) \mathbf{1}_p\|_2^2,$$

with sample version denoted $\widehat{R}(M)$. Consider all models that can be written as

$$M(X) = U \cdot D \cdot V(X)^\top$$

where U is an orthogonal $q \times r$ matrix, D is a positive diagonal matrix, and $V(X) = [v_{js}(X_j)]$ satisfies $\mathbb{E}(V^\top V) = I_r$. The population risk can be reexpressed as

$$\begin{aligned} R(M) &= \text{tr} \left\{ \begin{pmatrix} -I_q \\ DU^\top \end{pmatrix}^\top \mathbb{E} \left[\begin{pmatrix} Y \\ V(X)^\top \end{pmatrix} \begin{pmatrix} Y \\ V(X)^\top \end{pmatrix}^\top \right] \begin{pmatrix} -I_q \\ DU^\top \end{pmatrix} \right\} \\ &= \text{tr} \left\{ \begin{pmatrix} -I_q \\ DU^\top \end{pmatrix}^\top \begin{pmatrix} \Sigma_{YY} & \Sigma_{YV} \\ \Sigma_{YV}^\top & \Sigma_{VV} \end{pmatrix} \begin{pmatrix} -I_q \\ DU^\top \end{pmatrix} \right\} \end{aligned}$$

and similarly for the sample risk, with $\widehat{\Sigma}_n(V)$ replacing $\Sigma(V) := \text{Cov}((Y, V(X)^\top))$ above. The “uncontrollable” contribution to the risk, which does not depend on M , is $R_u = \text{tr}\{\Sigma_{YY}\}$. We can express the remaining “controllable” risk as

$$R_c(M) = R(M) - R_u = \text{tr} \left\{ \begin{pmatrix} -2I_q \\ DU^\top \end{pmatrix}^\top \Sigma(V) \begin{pmatrix} \mathbf{0}_q \\ DU^\top \end{pmatrix} \right\}.$$

Using the von Neumann trace inequality, $\text{tr}(AB) \leq \|A\|_p \|B\|_{p'}$ where $1/p + 1/p' = 1$,

$$\begin{aligned} R_c(M) - \widehat{R}_c(M) &\leq \left\| \begin{pmatrix} -2I_q \\ DU^\top \end{pmatrix}^\top (\Sigma(V) - \widehat{\Sigma}_n(V)) \right\|_{\text{sp}} \left\| \begin{pmatrix} \mathbf{0}_q \\ DU^\top \end{pmatrix} \right\|_* \\ &\leq \left\| \begin{pmatrix} -2I_q \\ DU^\top \end{pmatrix}^\top \right\|_{\text{sp}} \left\| \Sigma(V) - \widehat{\Sigma}_n(V) \right\|_{\text{sp}} \|D\|_* \\ &\leq C \max(2, \|D\|_{\text{sp}}) \left\| \Sigma(V) - \widehat{\Sigma}_n(V) \right\|_{\text{sp}} \|D\|_* \\ &\leq C \max\{2, \|D\|_*^2\} \left\| \Sigma(V) - \widehat{\Sigma}_n(V) \right\|_{\text{sp}} \end{aligned} \tag{7.1}$$

where here and in the following C is a generic constant. For the last factor in (7.1), it holds that

$$\sup_V \left\| \Sigma(V) - \widehat{\Sigma}_n(V) \right\|_{\text{sp}} \leq C \sup_V \sup_{w \in \mathcal{N}} w^\top (\Sigma(V) - \widehat{\Sigma}_n(V)) w$$

where \mathcal{N} is a $1/2$ -covering of the unit $(q+r)$ -sphere, which has size $|\mathcal{N}| \leq 6^{q+r} \leq 36^q$; compare [Vershynin \(2012, p. 665\)](#). We now assume that the functions $v_{sj}(x_j)$ are uniformly bounded from a Sobolev space of order two. Specifically, let $\{\psi_{jk} : k = 0, 1, \dots\}$ denote a uniformly bounded, orthonormal basis with respect to $L^2[0, 1]$, and assume that $v_{sj} \in \mathcal{H}_j$ where

$$\mathcal{H}_j = \left\{ f_j : f_j(x_j) = \sum_{k=0}^{\infty} a_{jk} \psi_{jk}(x_j), \quad \sum_{k=0}^{\infty} a_{jk}^2 k^4 \leq K^2 \right\}$$

for some $0 < K < \infty$. The L_∞ -covering number of \mathcal{H}_j satisfies $\log \mathcal{N}(\mathcal{H}_j, \epsilon) \leq K/\sqrt{\epsilon}$.

Suppose that $Y - \mathbb{E}(Y | X) = W$ is Gaussian and the true regression function $\mathbb{E}(Y | X)$ is bounded. Then the family of random variables $Z_{(V,w)} := \sqrt{n} \cdot w^\top (\Sigma(V) - \hat{\Sigma}_n(V))w$ is sub-Gaussian and sample continuous. It follows from a result of [Cesa-Bianchi and Lugosi \(1999\)](#) that

$$\mathbb{E} \left(\sup_V \sup_{w \in \mathcal{N}} w^\top (\Sigma(V) - \hat{\Sigma}_n(V))w \right) \leq \frac{C}{\sqrt{n}} \int_0^B \sqrt{q \log(36) + \log(pq) + \frac{K}{\sqrt{\epsilon}}} d\epsilon$$

for some constant B . Thus, by Markov's inequality we conclude that

$$\sup_V \left\| \Sigma(V) - \hat{\Sigma}_n(V) \right\|_{\text{sp}} = O_P \left(\sqrt{\frac{q + \log(pq)}{n}} \right), \quad (7.2)$$

when n tends to infinity and q and p possibly change with n . If

$$\|M\|_* = \|D\|_* = o \left(\frac{n}{(q + \log(pq))} \right)^{1/4},$$

then returning to (7.1), this gives us a bound on $R_c(M) - \hat{R}_c(M)$ that is $o_P(1)$. More precisely, for $\beta_n > 0$, define the class of matrices of functions

$$\mathcal{M}(\beta_n) = \{M : M(X) = UDV(X)^\top, \text{ with } \mathbb{E}(V^\top V) = I, v_{sj} \in \mathcal{H}_j, \|D\|_* \leq \beta_n\}. \quad (7.3)$$

Then, for a fitted matrix \widehat{M} chosen from $\mathcal{M}(\beta_n)$, writing $M_* = \arg \min_{M \in \mathcal{M}(\beta_n)} R(M)$, we have

$$\begin{aligned} R(\widehat{M}) - \inf_{M \in \mathcal{M}(\beta_n)} R(M) &= R(\widehat{M}) - \hat{R}(\widehat{M}) - (R(M_*) - \hat{R}(M_*)) + (\hat{R}(\widehat{M}) - \hat{R}(M_*)) \\ &\leq [R(\widehat{M}) - \hat{R}(\widehat{M})] - [R(M_*) - \hat{R}(M_*)]. \end{aligned}$$

Subtracting $R_u - \widehat{R}_u$ from each of the bracketed differences, we obtain that

$$\begin{aligned} R(\widehat{M}) - \inf_{M \in \mathcal{M}(\beta_n)} R(M) &\leq [R_c(\widehat{M}) - \hat{R}_c(\widehat{M})] - [R_c(M_*) - \hat{R}_c(M_*)] \\ &\leq 2 \sup_{M \in \mathcal{M}(\beta_n)} \{R_c(M) - \hat{R}_c(M)\} \\ &\stackrel{\text{by (7.1)}}{\leq} O_P \left(\|D\|_*^2 \left\| \Sigma(V) - \hat{\Sigma}_n(V) \right\|_{\text{sp}} \right). \end{aligned}$$

Now if

$$\beta_n = o \left(\frac{n}{q + \log(pq)} \right)^{1/4}, \quad (7.4)$$

then we may conclude from (7.2) that

$$R(\widehat{M}) - \inf_{M \in \mathcal{M}(\beta_n)} R(M) = o_P(1).$$

This proves the following result.

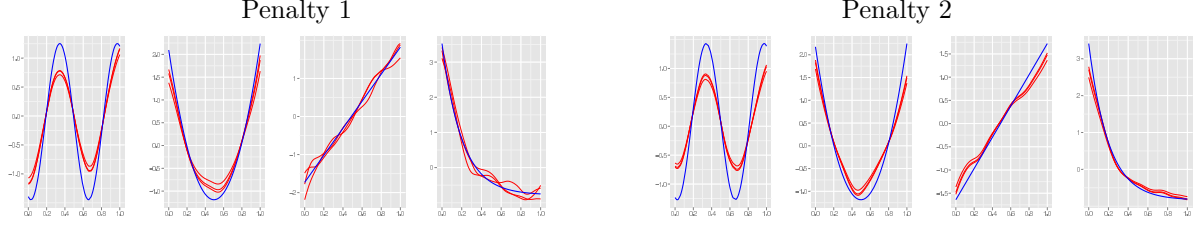


FIG 3. Example of the constrained-rank backfitting algorithm, with $q = 3$ and $p = 4$. The model is the same for each k ; that is, $m_j^k = m_j^{k'}$ for each j, k, k' . The left plots show the fits with penalties $\lambda_j \|M_j\|_*$, with $\lambda = (3, 3, 0, 0)$. The right plots use penalty 2, and the solution has rank 1.

Proposition 7.1. Let \widehat{M} minimize the empirical risk $\frac{1}{n} \sum_i \|Y_i - \sum_j M_j(X_{ij})\|_2^2$ over the class $\mathcal{M}(\beta_n)$. Suppose that $Y - \mathbb{E}(Y | X)$ is Gaussian, the true regression function $\mathbb{E}(Y | X)$ is bounded, and β_n satisfies (7.4) as $n \rightarrow \infty$. Then it holds that

$$R(\widehat{M}) - \inf_{M \in \mathcal{M}(\beta_n)} R(M) \xrightarrow{P} 0.$$

8. Examples

8.1. Synthetic Data Example

Figure 3 shows an example of the backfitting algorithms for penalties 1 and 2. For this example $q = 3$ and $p = 4$. The model is the same for each k ; that is, $m_j^k = m_j^{k'}$ for each j, k, k' . The true regression functions are $m_1^k(x) = \sin(2x)$, $m_2^k(x) = x^2 - c_2$, $m_3^k(x) = x$, and $m_4^k(x) = e^{-x} - c_4$, where c_2 and c_4 are centering constants. The left group of plots in Figure 3 shows the fits obtained with penalties $\lambda_j \|M_j\|_*$ with regularization parameters $\lambda = (\lambda_1, \lambda_2, \lambda_3, \lambda_4) = (3, 3, 0, 0)$, so that the third and fourth functions are not regularized. The plots show how the first and second function estimates are identical up to multiplicative scaling for each $k = 1, 2, 3$ (\mathbb{M}_1 and \mathbb{M}_2 have rank one), while the third and fourth function estimates vary (\mathbb{M}_3 and \mathbb{M}_4 have rank three). The fits are made with local linear smoothing. The right set of plots corresponds to penalty 2, and the solution has rank 1. We intentionally use a bandwidth that is too small, to better show the differences with and without regularization. The true regression functions m_j^k are shown in blue in the plots; the fitted functions are in red, with the fits for a given j superimposed on the same plot. The sample size is $n = 150$, and the noise variance is $\sigma^2 = 1$.

8.2. Gene Expression Data

To further illustrate the proposed nonparametric reduced rank regression techniques, we consider data on gene expression in *E. coli* from the “DREAM 5 Network Inference Challenge”¹

¹<http://wiki.c2b2.columbia.edu/dream/index.php/D5c4>

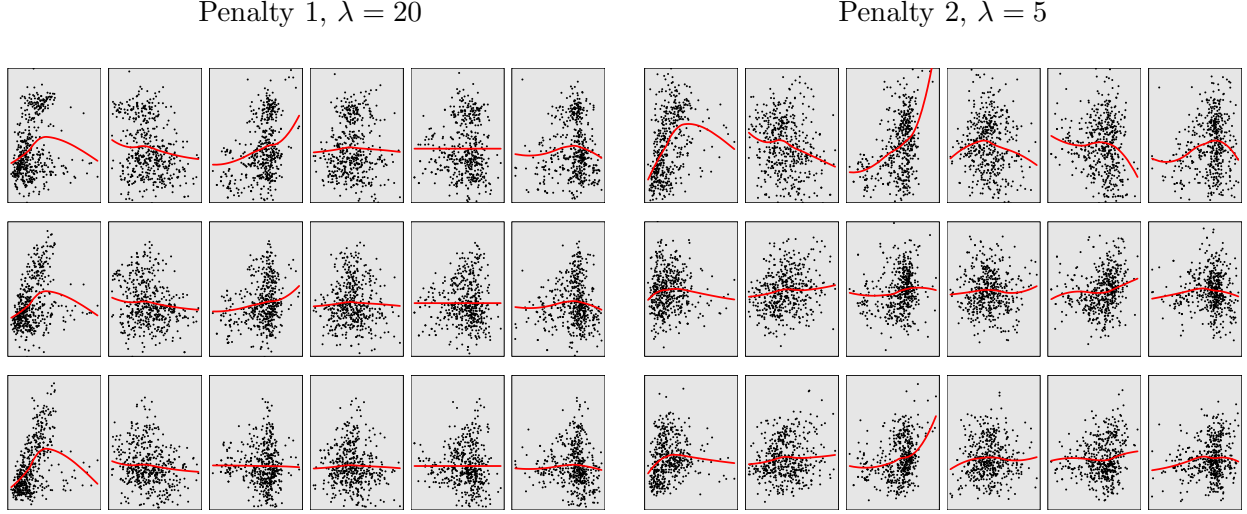


FIG 4. Fits on the gene expression data. Left: Penalty 1 with large tuning parameter. Right: Penalty 2 tuned via 10-fold cross-validation. Plotted points are residuals holding out the given predictor.

(Marbach *et al.*, 2012). In this challenge genes were classified as transcription factors (TFs) or target genes (TGs). Transcription factors regulate the target genes, as well as other TFs.

We focus on predicting the expression levels Y for a particular set of $q = 27$ TGs, using the expression levels X for $p = 6$ TFs. Our motivation for analyzing these 33 genes is that, according to the gold standard gene regulatory network used for the DREAM 5 challenge, the 6 TFs form the parent set common to two additional TFs, which have the 27 TGs as their child nodes. In fact, the two intermediate nodes d-separate the 6 TFs and the 27 TGs in a Bayesian network interpretation of this gold standard. This means that if we treat the gold standard as a causal network, then up to noise, the functional relationship between X and Y is given by the composition of a map $g : \mathbb{R}^6 \rightarrow \mathbb{R}^2$ and a map $h : \mathbb{R}^2 \rightarrow \mathbb{R}^{27}$. If g and h are both linear, their composition $h \circ g$ is a linear map of rank at most 2. As observed in Section 2, such a reduced rank linear model is a special case of an additive model with reduced rank in the sense of penalty 2. More generally, if g is an additive function and h is linear, then $h \circ g$ has rank at most 2 in the sense of penalty 2. Higher rank can in principle occur under functional composition, since even a univariate additive map $h : \mathbb{R} \rightarrow \mathbb{R}^q$ may have rank up to q under our penalties (penalty 1 and 2 coincide for univariate maps).

The backfitting algorithm of Figure 1 with penalty 1 and a rather aggressive choice of the tuning parameter λ produces the estimates shown in Figure 4, for which we have selected three of the 27 TGs. Under such strong regularization, the 5th column of functions is rank zero and, thus, identically zero. The remaining columns have rank one; the estimated fitted values are scalar multiples of one another. We also see that scalings can be different for different columns. The third plot in the third row shows a slightly negative slope, indicating a negative scaling for this particular estimate. The remaining functions in this row are oriented similarly to the other rows, indicating the same, positive scaling. This property characterizes

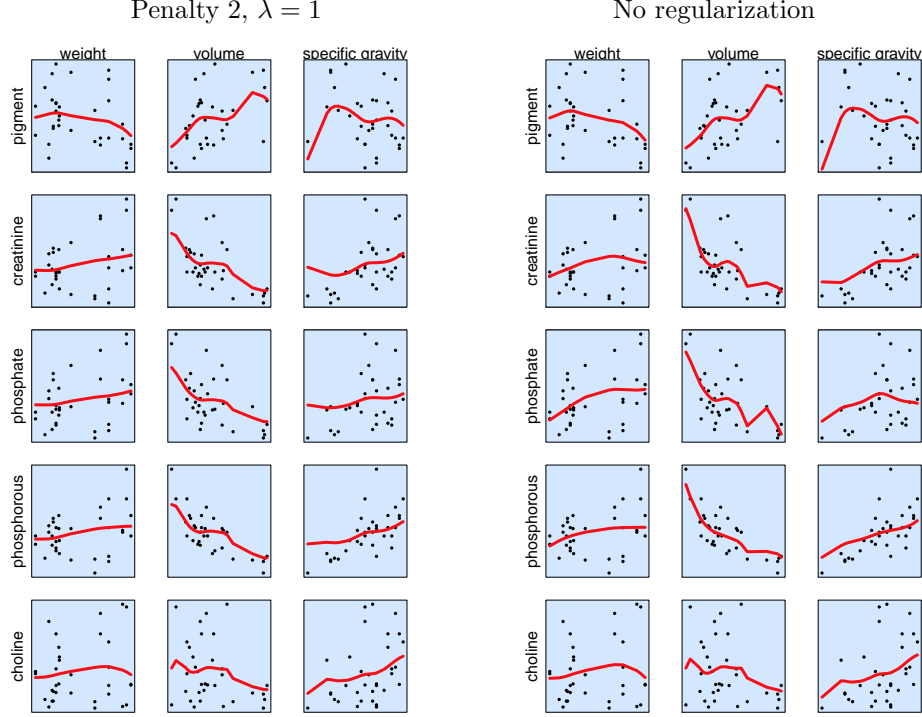


FIG 5. Fits on the biochemistry data with penalty 2 (left), with $\lambda = 1$. The solution has rank three, with regression functions for creatinine, phosphate and phosphorous identical up to scaling. The fits on the right have no regularization, and are full rank.

the difference between penalties 1 and 2; in an application of penalty 2, the scalings would have been the same across all functions in a given row.

Next, we illustrate a higher-rank solution for penalty 2. Choosing the regularization parameter λ by ten-fold cross-validation gives a fit of rank 5, considerably lower than 27, the maximum possible rank. Figure 4 shows a selection of three coordinates of the fitted functions. Under rank five, each row of functions is a linear combination of up to five other, linearly independent rows. We remark that the use of cross-validation generally produces somewhat more complex models than is necessary to capture an underlying low-rank data-generating mechanism. Hence, if the causal relationships for these data were indeed additive and low rank, then the true low rank might well be smaller than five.

8.3. Biochemistry Example

Here we analyze the same biochemical data of Smith et al. (1962) that was used by Yuan et al. (2007). The data contain chemical measurements for 33 individual samples of men's urine specimens. The $q = 5$ response variables are pigment creatinine, and the concentrations (in mg/ml) of phosphate, phosphorous, creatinine and choline. The $p = 3$ covariates are the weight of the subject, and volume and specific gravity of the specimen. Yuan et al. (2007) form a linear model where the coefficients in a spline basis are regularized. Their

plots suggest some boundary effects, due to the choice of basis. Here we use our backfitting algorithm with local linear smoothing. No explicit basis is used. Figure 5 shows the result of using regularization $\lambda \|M_{1:p}\|_*$ with $\lambda = 1$ (penalty 2), and bandwidth $h = .3$ for all variables. The regularized solution has rank 3, where the response variables for phosphorous, phosphate, and creatinine concentrations are scaled versions of each other. The plots on the right show fits with no regularization ($\lambda = 0$).

9. Summary

This paper introduced two penalties that induce reduced rank fits in multivariate additive nonparametric regression. Under linearity, the penalties specialize to group lasso and nuclear norm penalties for classical reduced rank regression. Examining the subdifferentials of each of these penalties, we developed backfitting algorithms for the two resulting optimization problems that are based on soft-thresholding of singular values of smoothed residual matrices. The algorithms were demonstrated on a gene expression data set and a biochemical data set that have low-rank structure. We also provided a persistence analysis that shows error tending to zero under a scaling assumption on the sample size n and the dimensions q and p .

Acknowledgements

Research supported in part by NSF grants IIS-1116730, DMS-0746265, and DMS-1203762, AFOSR grant FA9550-09-1-0373, ONR grant N000141210762, and an Alfred P. Sloan Fellowship.

References

- CESA-BIANCHI, N. and LUGOSI, G. (1999). On prediction of individual sequences. *The Annals of Statistics* **27** 1865–1894.
- DINUZZO, F. and FUKUMIZU, K. (2012). Learning low-rank output kernels. *Journal of Machine Learning Research* **20** 181–196. Asian Conference on Machine Learning proceedings.
- FAZEL, M. (2002). Matrix rank minimization with applications. Doctoral Dissertation, Electrical Engineering Department, Stanford University.
- FOYGEL, R., HORRELL, M., DRTON, M. and LAFFERTY, J. (2012). Nonparametric Reduced Rank Regression. In *Advances in Neural Information Processing Systems 25* (P. Bartlett, F. Pereira, C. Burges, L. Bottou and K. Weinberger, eds.) 1637–1645.
- MARBACH, D., COSTELLO, J. C., KÜFFNER, R., VEGA, N., PRILL, R. J., CAMACHO, D. M., ALLISON, K. R., THE DREAM5 CONSORTIUM, KELLIS, M., COLLINS, J. J. and STOLOVITZKY, G. (2012). Wisdom of crowds for robust gene network inference. *Nature Methods* **9** 796–804.
- NEGAHBAN, S. and WAINWRIGHT, M. J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Annals of Statistics* **39** 1069–1097.

- RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2012). Minimax-optimal rates for sparse additive models over kernel classes via convex programming. *J. Mach. Learn. Res.* **13** 389–427.
- RAVIKUMAR, P., LAFFERTY, J., LIU, H. and WASSERMAN, L. (2009). Sparse additive models. *Journal of the Royal Statistical Society, Series B, Methodological* **71** 1009–1030.
- RECHT, B., FAZEL, M. and PARRILO, P. A. (2010). Guaranteed minimum rank solutions to linear matrix equations via nuclear norm minimization. *SIAM Review* **52** 471–501.
- VERSHYNIN, R. (2012). How close is the sample covariance matrix to the actual covariance matrix? *J. Theoret. Probab.* **25** 655–686.
- WATSON, G. A. (1992). Characterization of the subdifferential of some matrix norms. *Linear Algebra and Applications* **170** 1039–1053.
- YUAN, M. and LIN, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68** 49–67.
- YUAN, M., EKICI, A., LU, Z. and MONTEIRO, R. (2007). Dimension reduction and coefficient estimation in multivariate linear regression. *J. R. Statist. Soc. B* **69** 329–346.